

Decision Trees

Introduction

What is Segmentation?

What is Segmentation?

- Imagine a scenario where we want to run a SMS marketing campaign to attract more customers in the next quarter
 - Some customers like to see high discount
 - Some customers want to see a large collection of items
 - Some customers are fans of particular brands
 - Some customers are Male some are Female
- Divide them based on their demographics, buying patterns and profile related attributes

What is Segmentation?

- One size doesn't fit all
- Divide the population in such a way that
 - Customers inside a group are homogeneous
 - Customers across groups are heterogeneous
- Is there any statistical way of dividing them correctly based on the data

Segmentation Business Problem

The Business Problem

Old Data

Gender	Marital Status	Ordered the product
M	Married	No
F	Unmarried	Yes
M	Married	No
M	Married	No
M	Married	No
M	Married	No
F	Unmarried	Yes
M	Unmarried	Yes
F	Married	No
M	Married	No
F	Married	No
M	Unmarried	No
F	Married	No
F	Unmarried	Yes

New Data

Gender	Marital Status	Product order
M	Married	??
F	Unmarried	??

The Business Problem

Old Data

Sr No	Gender	Marital Status	Ordered the product
1	M	Married	No
2	F	Unmarried	Yes
3	M	Married	No
4	M	Married	No
5	M	Married	No
6	M	Married	No
7	F	Unmarried	Yes
8	M	Unmarried	Yes
9	F	Married	No
10	M	Married	No
11	F	Married	No
12	M	Unmarried	No
13	F	Married	No
14	F	Unmarried	Yes

New Data

Gender	Marital Status	Product order
M	Married	??
F	Unmarried	??

The Decision Tree Philosophy

The Data

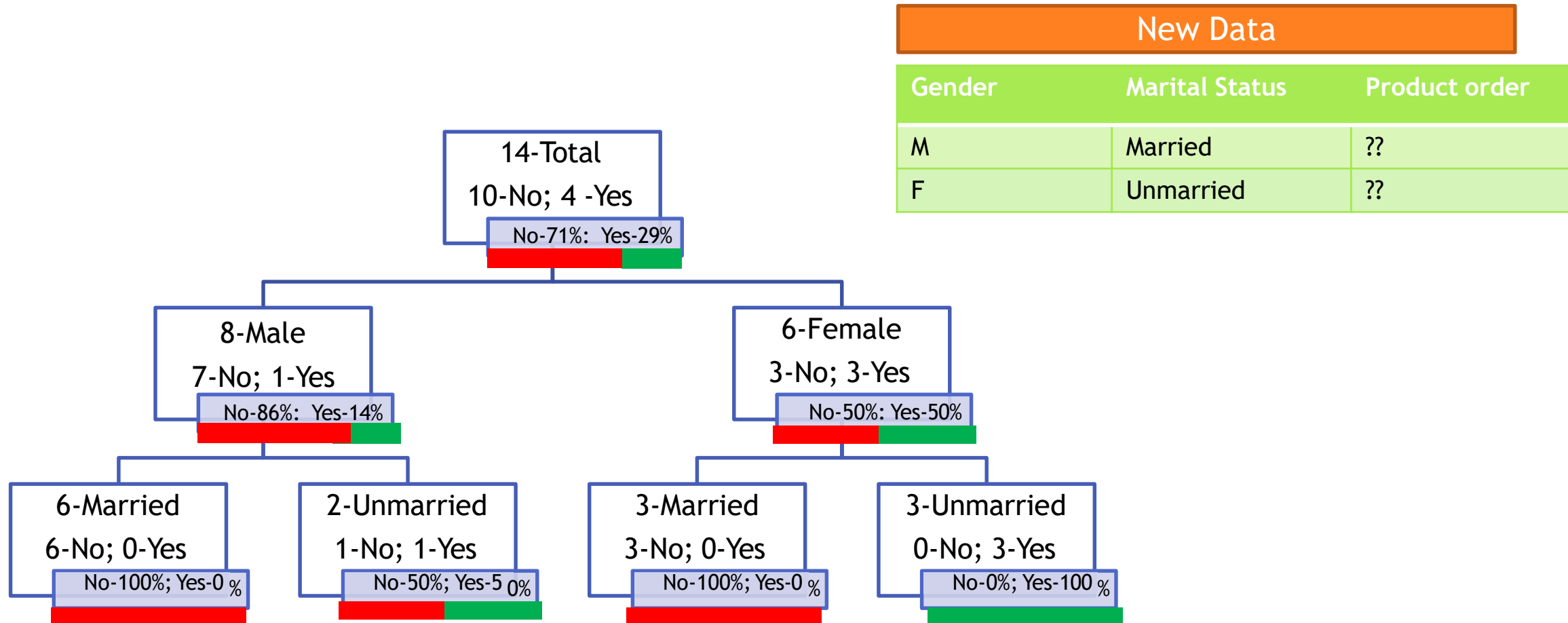
Old Data

Sr No	Gender	Marital Status	Ordered the product
1	M	Married	No
2	F	Unmarried	Yes
3	M	Married	No
4	M	Married	No
5	M	Married	No
6	M	Married	No
7	F	Unmarried	Yes
8	M	Unmarried	Yes
9	F	Married	No
10	M	Married	No
11	F	Married	No
12	M	Unmarried	No
13	F	Married	No
14	F	Unmarried	Yes

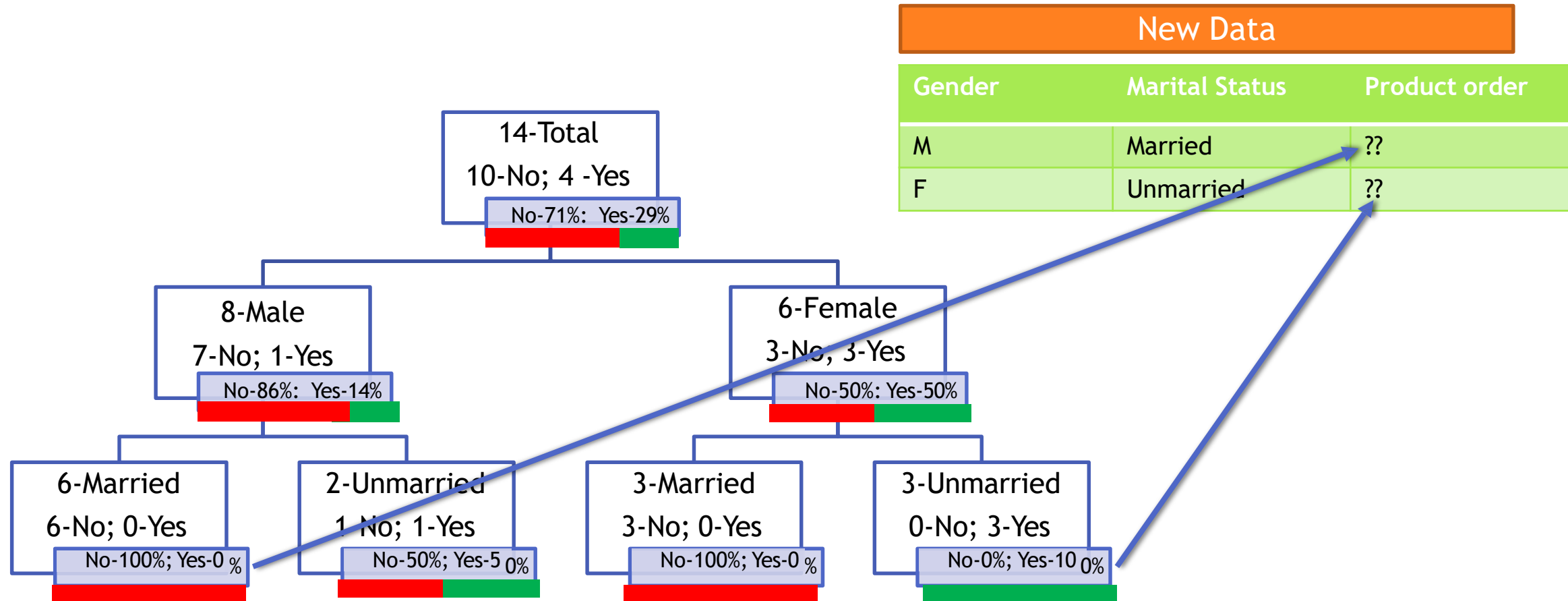
New Data

Gender	Marital Status	Product order
M	Married	??
F	Unmarried	??

Re-Arranging the data



Re-Arranging the data

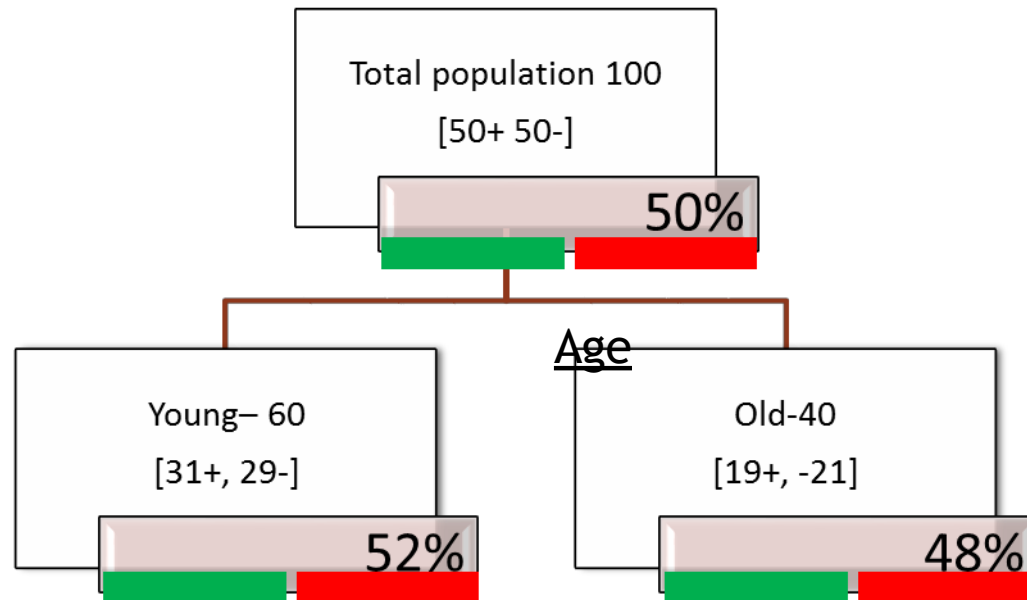


The Decision Tree Approach

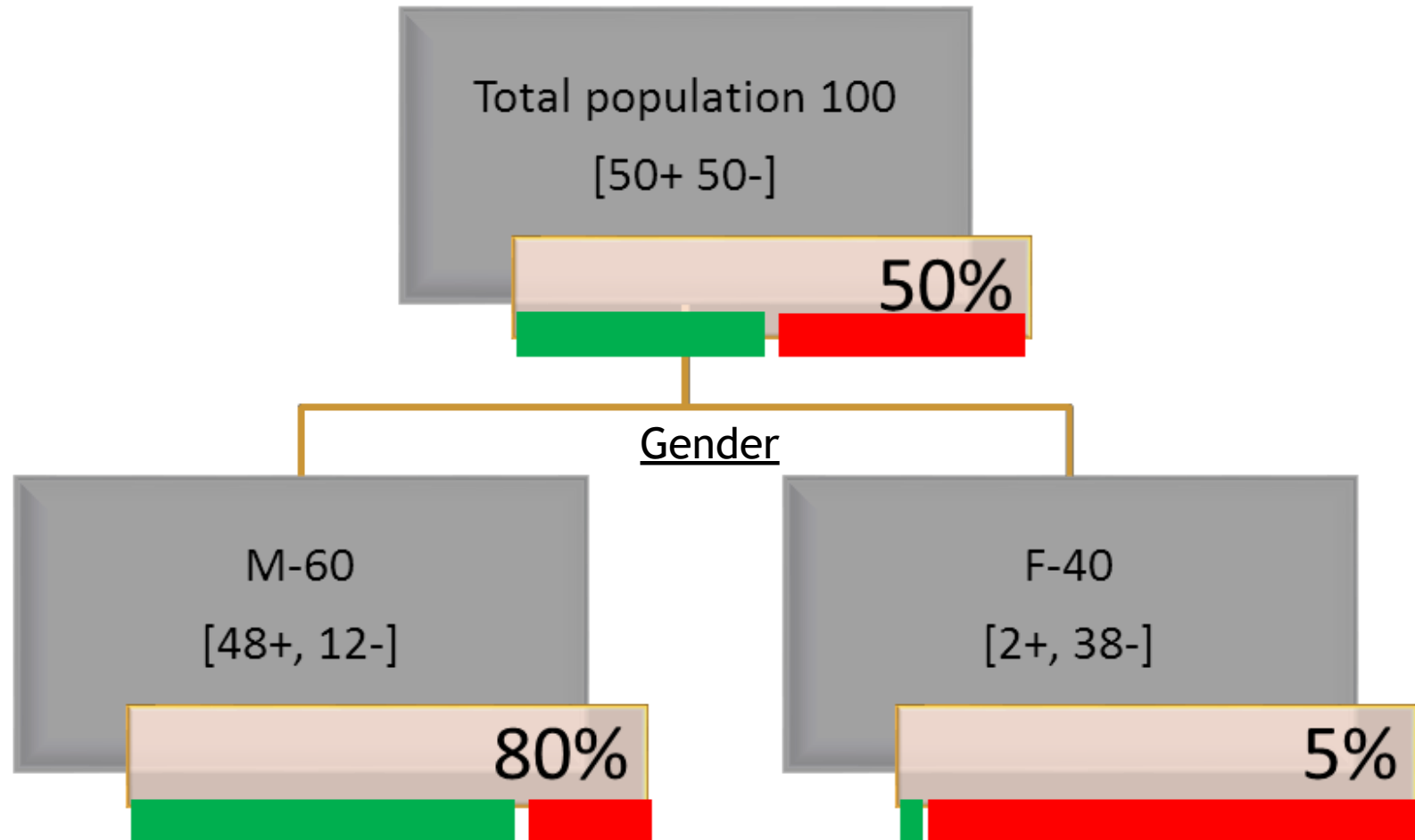
The Decision Tree Approach

- The aim is to decide the whole population or the data set into segments
- The segmentation need to be useful for business decision making.
- If one class is really dominating in a segments
 - Then it will be easy for us to classify the unknown items
 - Then its very easy for applying business strategy
- For example:
 - It takes no great skill to say that the customers have 50% chance to buy and 50% chance to not buy.
 - A good splitting criterion segments the customers with 90% -10% buying probability, say Gender=“Female” customers have 5% buying probability and 95% not buying

Example Sales Segmentation Based on Age



Example Sales Segmentation Based on Gender



Main questions

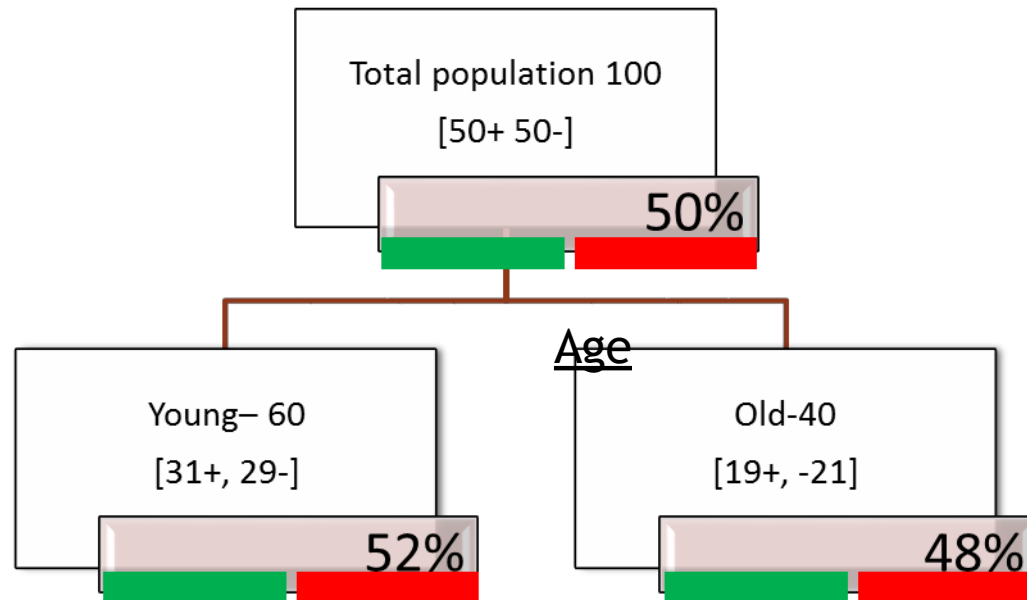
- Ok we are looking for pure segments
- Dataset has many attributes
- Which is the right attribute for pure segmentation?
- Can we start with any attribute?
- Which attribute to start? - The best separating attribute
- Customer Age can impact the sales, gender can impact sales , customer place and demographics can impact the sales. How to identify the best attribute and the split?

The Splitting Criterion

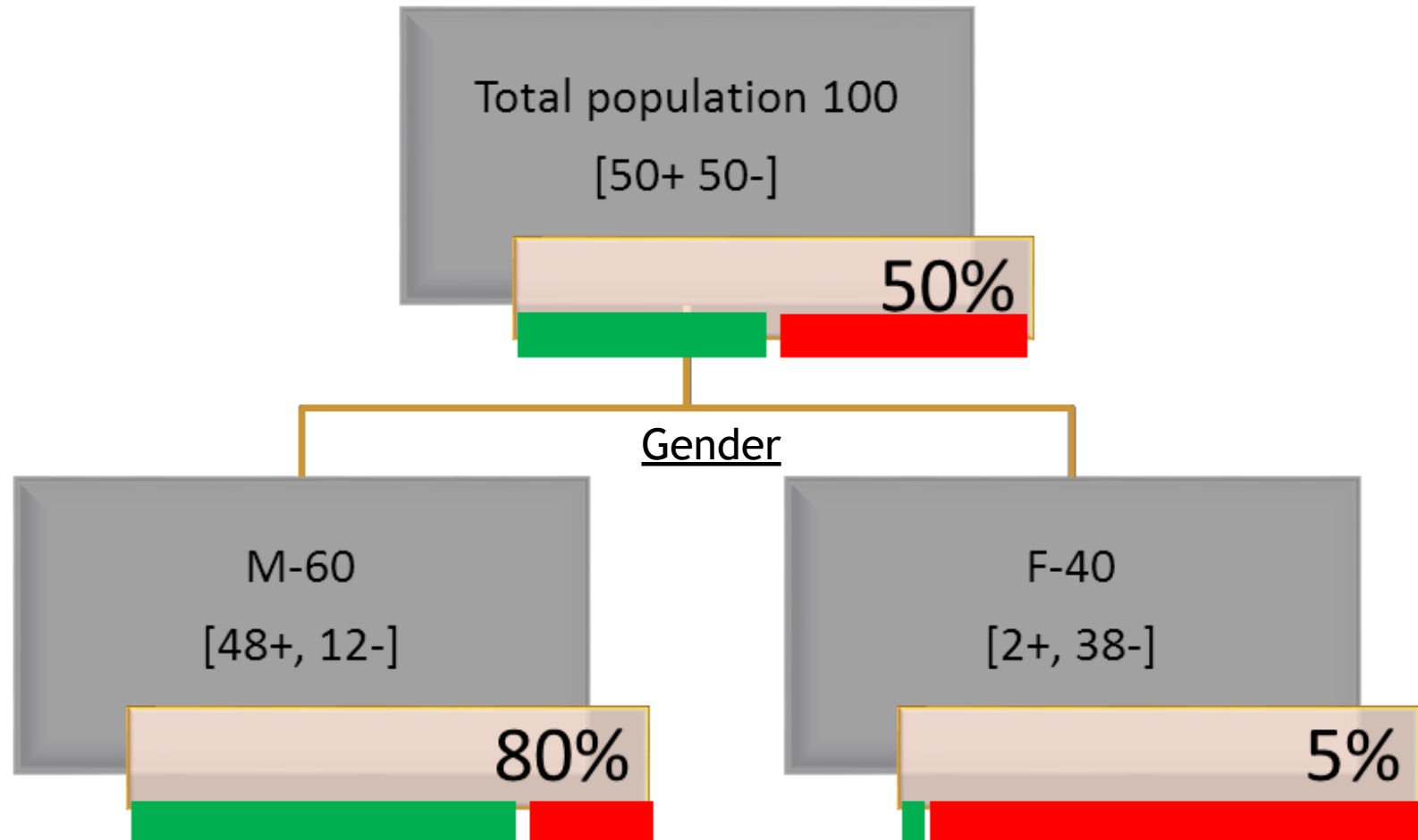
The Splitting Criterion

- The best split is
 - The split does the best job of separating the data into groups
 - Where a single class (either 0 or 1) predominates in each group

Example Sales Segmentation Based on Age

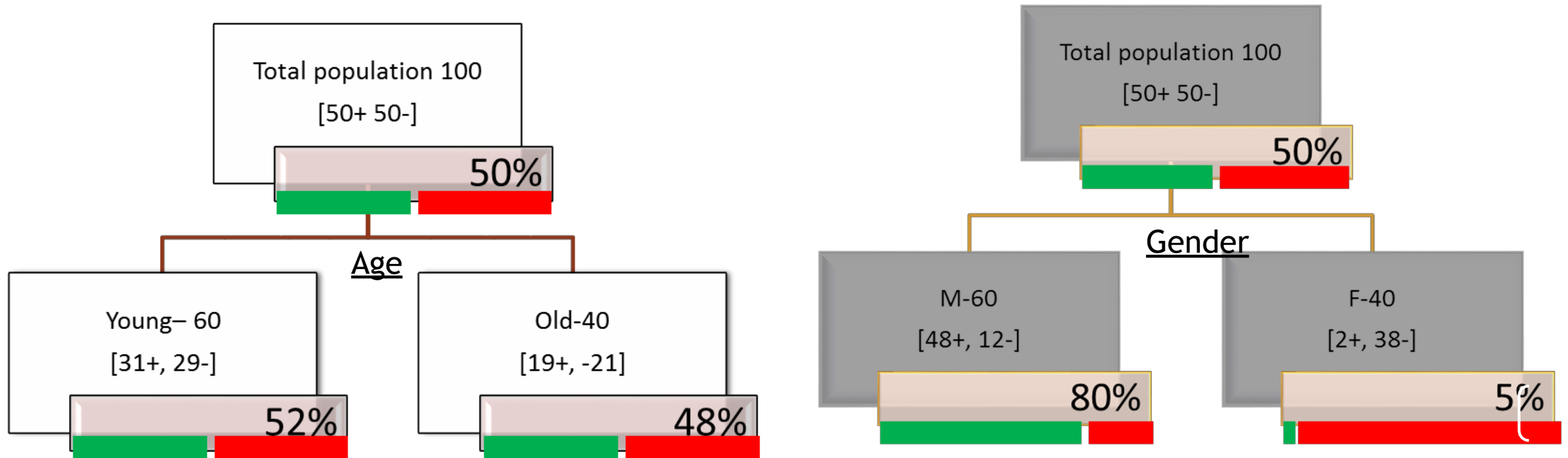


Example Sales Segmentation Based on Gender



Impurity (Diversity) Measures:

- We are looking for a impurity or diversity measure that will give high score for this Age variable (high impurity while segmenting), Low score for Gender variable (Low impurity while segmenting)

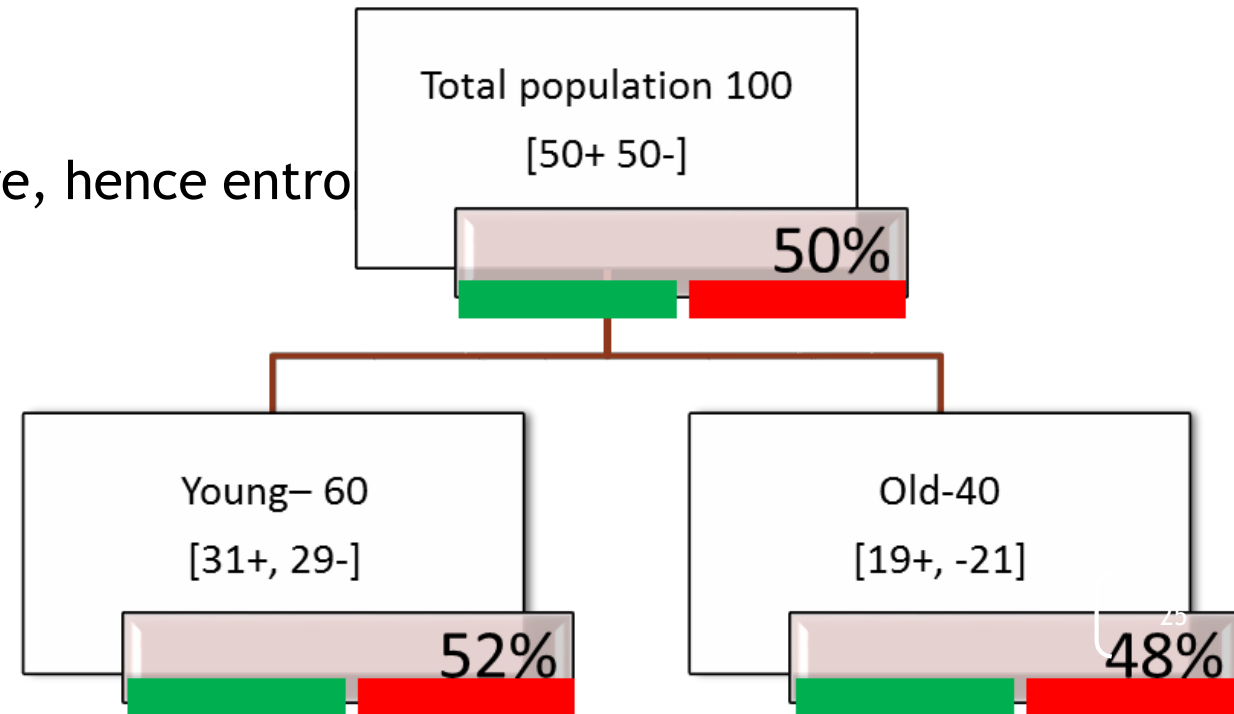


Impurity (Diversity) Measures:

- **Entropy:** Characterizes the impurity/diversity of segment
- Measure of uncertainty/Impurity
- Entropy measures the information amount in a message
- S is a segment of training examples, p_+ is the proportion of positive examples, p_- is the proportion of negative examples
- $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
 - Where p_+ is the probability of positive class and p_- is the probability of negative class
- Entropy is highest when the split has p of 0.5.
- Entropy is least when the split is pure .ie p of 1

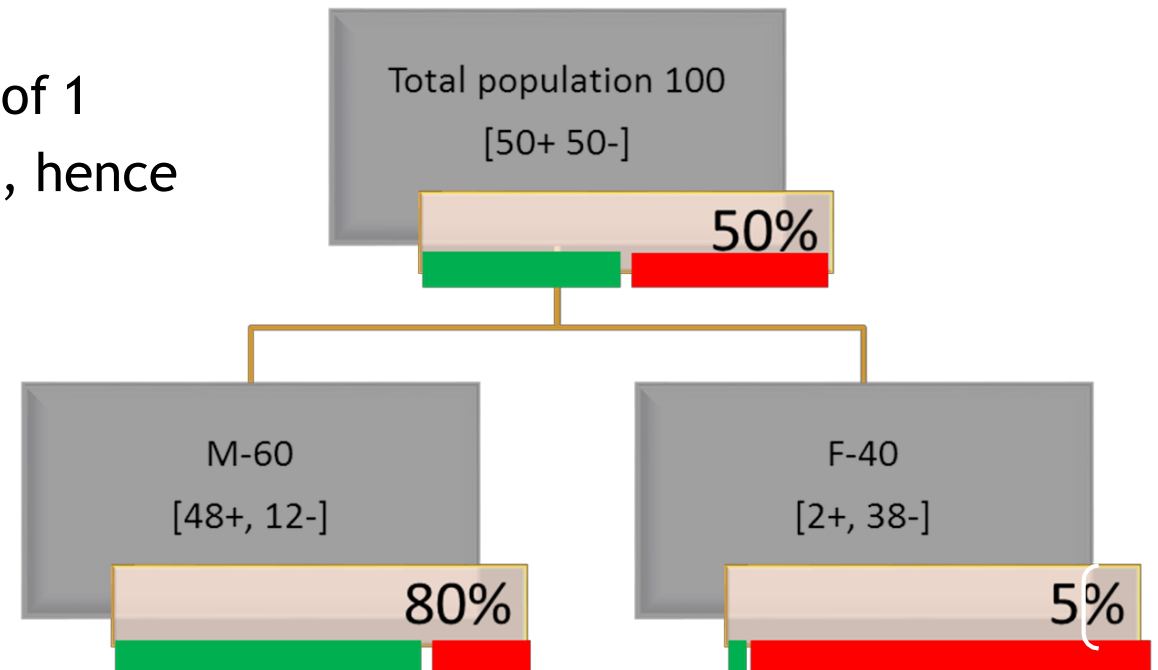
Entropy is highest when the split has p of 0.5

- S is a segment of training examples, p_+ is the proportion of positive examples, p_- is the proportion of negative examples
- $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
- Entropy is highest when the split has p of 0.5
- 50-50 class ratio in a segment is really impure, hence entropy
 - $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
 - $\text{Entropy}(S) = -0.5 \log_2(0.5) - 0.5 \log_2(0.5)$
 - $\text{Entropy}(S) = 1$



Entropy is least when the split is pure .ie p of 1

- S is a segment of training examples, p_+ is the proportion of positive examples, p_- is the proportion of negative examples
- $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
 - Entropy is least when the split is pure .ie p of 1
 - 100-0 class ratio in a segment is really pure, hence entropy is low
 - $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
 - $\text{Entropy}(S) = -1 \cdot \log_2(1) - 0 \cdot \log_2(0)$
 - $\text{Entropy}(S) = 0$



The less the entropy, the better the split

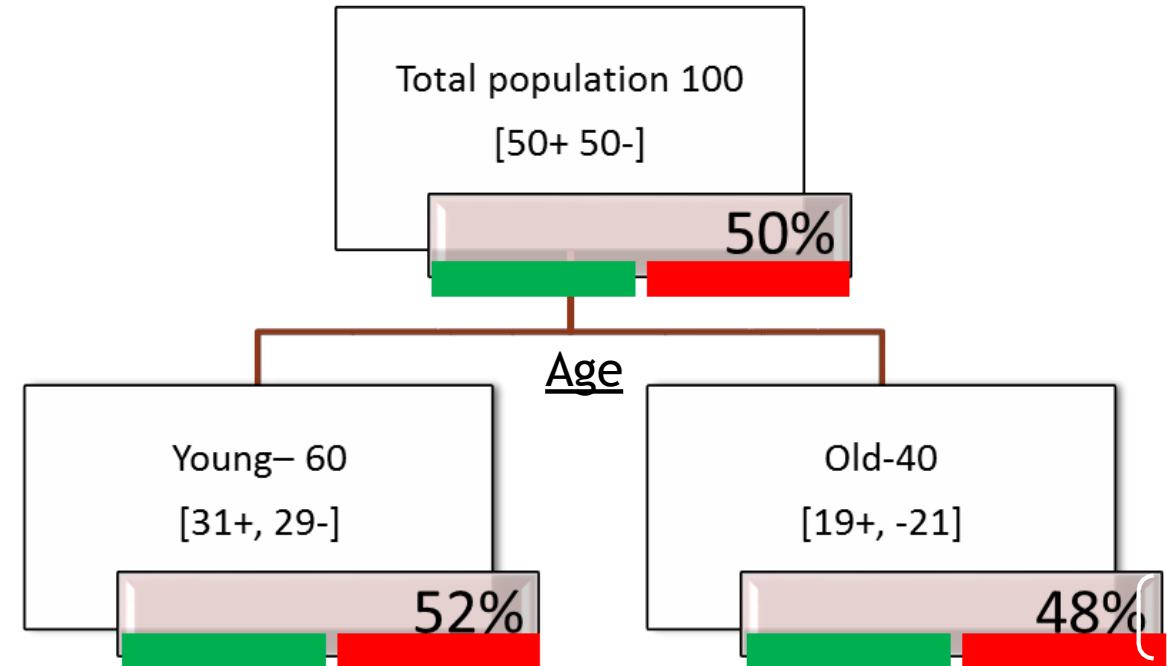
- The less the entropy, the better the split
- Entropy is formulated in such a way that, its value will be high for impure segments

Entropy Calculation - Example

Entropy Calculation

- Entropy at root
- Total population at root 100 [50+,50-]
- $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
- $-0.5 \log_2 (0.5) - 0.5 \log_2 (0.5)$
- $-(0.5)*(-1) - (0.5)*(-1)$
- 1
- 100% Impurity at root

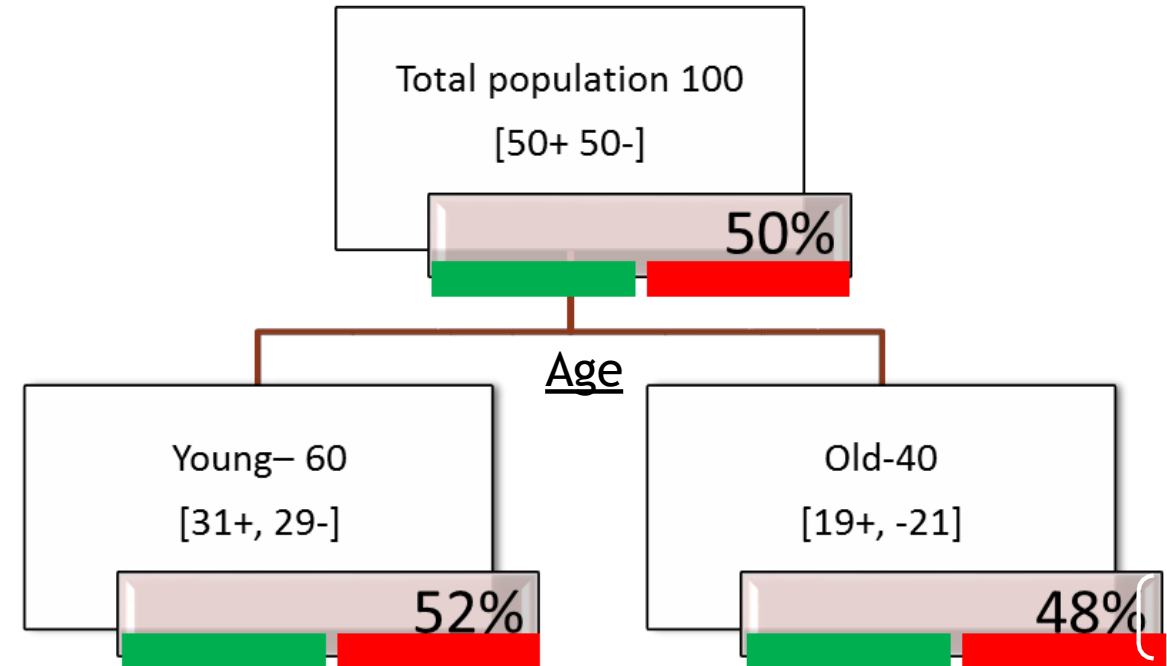
$$\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$



Entropy Calculation

- Gender Splits the population into two segments
- Segment-1 : Age="Young"
- Segment-2: Age="Old"
- Entropy at segment-1
 - Age="Young" segment has 60 records [31+,29-]
 - $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
 - $-31/60 \log_2 31/60 - 29/60 \log_2 29/60$
 - $(-31/60) \cdot \log(31/60, 2) - (29/60) \cdot \log(29/60, 2)$
 - 0.9991984 (99% Impurity in this segment)

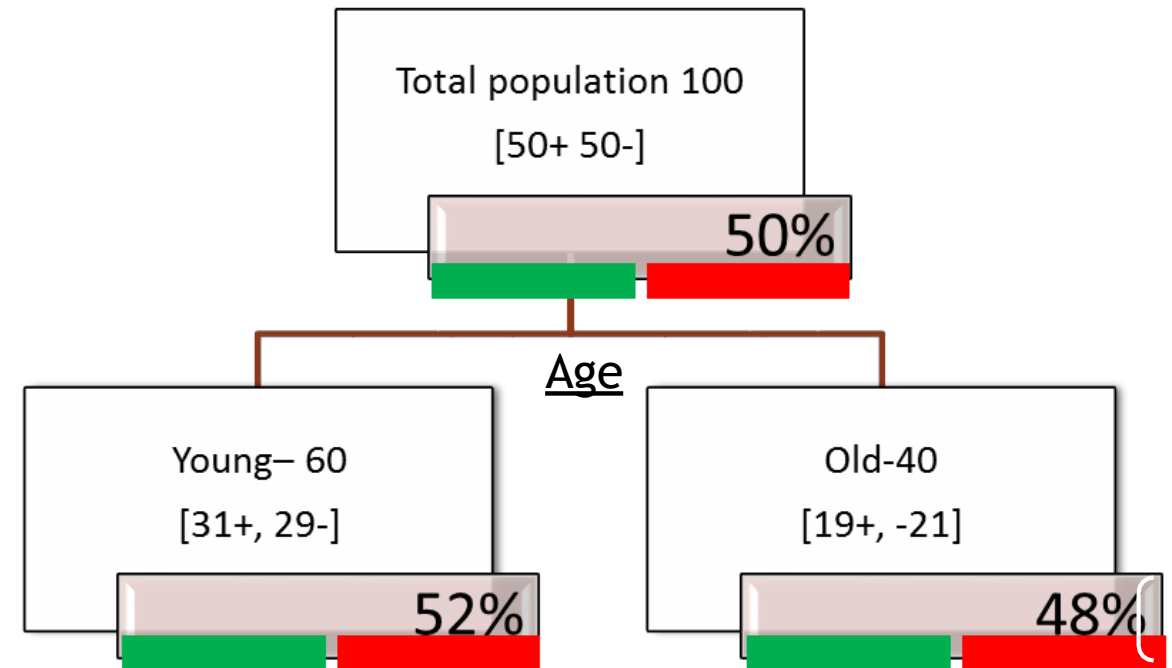
$$\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$



Entropy Calculation

- Gender Splits the population into two segments
- Segment-1 : Age="Young"
- Segment-2: Age="Old"
- Entropy at segment-2
 - Age="Old" segment has 40 records [19+,21-]
 - $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
 - $-19/40 \log_2 19/40 - 21/40 \log_2 21/40$
 - $(-19/40) * \log(19/40, 2) - (21/40) * \log(21/40, 2)$
 - 0.9981959 (99% Impurity in this segment too)

$$\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$



LAB: Decision Tree Building

LAB: Decision Tree Building

- Data:Ecom_Cust_Relationship_Management/Ecom_Cust_Survey.csv
- How many customers have participated in the survey?
- Overall most of the customers are satisfied or dis-satisfied?
- Can you segment the data and find the concentrated satisfied and dis-satisfied customer segments ?
- What are the major characteristics of satisfied customers?
- What are the major characteristics of dis-satisfied customers?

Spyder (Python 3.5)

File Edit Search Source Run Debug Consoles Projects Tools View Help

C:\Users\StatInfer

Editor - D:\Dropbox\B. Video Proj Files\Python ML\Session 3. Decision Tree\Active Presenter files\3.Decision_Trees_V3.py

3.Decision_Trees_V3.py*

```
64 #We will need to install graphviz tool in our system and set the path in environment variables.
65 #Visit http://www.graphviz.org/Download..php and find the optimal version for the computer.
66 #Get the path for gvedit.exe in install directory(for me it was "C:\Program Files (x86)\Graphviz2.38\
67 #goto start->computer->system properties->advanced settings->environment variables and add the path.
68 #We will need python package pydotplus(for older python versions pydot)
69 #use this command in your anaconda prompt: conda install -c conda-forge pydotplus
70 #if an error regarding version occurs while installing the package go to https://anaconda.org/search/
71 #this link will show the channel name of the suitable version suitable.
72 #and we can use use : conda install -c <channel name here> pydotplus
73
74
75 import pydotplus
76 dot_data = StringIO()
77 tree.export_graphviz(clf,
78                     out_file = dot_data,
79                     feature_names = features,
80                     filled=True, rounded=True,
81                     impurity=False)
82
83 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
84 Image(graph.create_png())
85
86 #by looking at the plot we can answer two questions:
87 #Q.5 What are the major characteristics of satisfied customers?
88
89 #Major_characteristics= Order.Quantity<40 & Age<30 / Order.Quantity >=40
90
91 #Q 6. What are the major characteristics of dis-satisfied customers?
92 #Major_characteristics= Order.Quantity<40 & Age>=30
93
94 #IAB : Tree Validation
95 #####
96 #####Tree Validation
97 #Tree Validation
98 predict1 = clf.predict(X)
99 predict1
100
101 from sklearn.metrics import confusion_matrix ###for using confusion matrix###
102 cm = confusion_matrix(y, predict1)
103 print (cm)
104
105 total = sum(sum(cm))
```

Variable explorer

Name	Type	Size	Value
X	DataFrame	(11805, 4)	Column names: Age, Order_Quantity, Customer_Type, Improvement_Area
df	DataFrame	(11805, 7)	Column names: Cust_num, Region, Age, Order_Quantity, Customer_Type, Improvement_Area, Overall_Satisfaction
features	list	4	['Age', 'Order_Quantity', 'Customer_Type', 'Improvement_Area']
satisfied	Int64	1	class 'numpy.int64'
y	Series	(11805,)	class 'pandas.core.series.Series'

Variable explorer File explorer Help

Python console

Console 1/A

```
....:
....: graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
....: Image(graph.create_png())
....:
Out[44]:
```

ROOT NODE

```
graph TD
    Root["Order_Quantity <= 40.5  
samples = 11805  
value = [6408, 5397]"]
    True["Age <= 29.5  
samples = 7397  
value = [6374, 1023]"]
    False["Age <= 20.5  
samples = 4408  
value = [34, 4374]"]
    L1["samples = 379  
value = [4, 375]"]
    L2["samples = 7018  
value = [6370, 648]"]
    L3["samples = 475  
value = [6, 469]"]
    L4["samples = 3933  
value = [28, 3905]"]

    Root -- True --> True
    Root -- False --> False
    True --> L1
    True --> L2
    False --> L3
    False --> L4
```

In [45]:

```

Spyder (Python 3.5)
File Edit Search Source Run Debug Consoles Projects Tools View Help
C:\Users\StatInfer
Editor - D:\Dropbox\0. Video Proj Files\Python ML\Session 3. Decision Tree\Active Presenter files\3.Decision_Trees_V3.py
3.Decision_Trees_V3.py*
1 #Import Data
2 import pandas as pd
3
4 ##Ecom_Cust_Survey = pd.read_csv('...',header = 0)
5 df = pd.read_csv('D:\\Google Drive\\Training\\Datasets\\Ecom_Cust_Relationship_Management\\Ecom_Cust_
6
7 df.dropna(inplace=True) # to remove all the missing values rows..
8 #Q 1. How many customers have participated in the survey?
9 df.info()
10 #ANS: 11805
11
12 #Q.2 Overall most of the customers are satisfied or dis-satisfied?
13 #total number of customers
14 df.shape
15 df.head()
16
17
18 #number of satisfied customers
19 satisfied = df['Overall_Satisfaction'].map( {'Dis Satisfied': 0, 'Satisfied': 1} ).astype(int).sum()
20 satisfied
21 #number of dis-satisfied customers
22 df.shape[0]-satisfied
23
24
25 #Q 3. Can you segment the data and find the concentrated satisfied and dis-satisfied customer segments
26 #solution:
27 # We will create a tree model in python using the sci-kit module
28 # before that we will need to convert most of the feature data into numerical or hash values as scikit
29 # Welcome to variable transformation
30
31 df['Region'] = df['Region'].map( {'EAST': 1, 'WEST': 2, 'NORTH': 3, 'SOUTH':4} ).astype(int)
32 df['Customer_Type'] = df['Customer_Type'].map({'Prime': 1, 'Non_Prime': 0}).astype(int)
33
34 #We will also need to change the column names, as '.' and spaces are part of many basic functions in py
35 df.rename(columns={'Order Quantity':'Order_Quantity', 'Improvement Area':'Improvement_Area'}, inplace
36 df['Improvement_Area'] = df['Improvement_Area'].map({'Website UI':1, 'Packing & Shipping':2, 'Product
37 df['Overall_Satisfaction'] = df['Overall_Satisfaction'].map( {'Dis Satisfied': 0, 'Satisfied': 1} ).as
38
39 #Need the library to create the tree
40 from sklearn.tree import DecisionTreeClassifier
41
42 #Defining Features and Labels

```

Variable explorer

Name	Type	Size	Value
X	DataFrame	(11805, 4)	Column names: Age, Order_Quantity, Customer_Type, Improvement_Area
df	DataFrame	(11805, 7)	Column names: Cust_num, Region, Age, Order_Quantity, Customer_Type, Improvement_Area, Overall_Satisfaction
features	list	4	['Age', 'Order_Quantity', 'Customer_Type', 'Improvement_Area']
satisfied	int64	1	class 'numpy.int64'
y	Series	(11805,)	class 'pandas.core.series.Series'

Variable explorer File explorer Help

IPython console

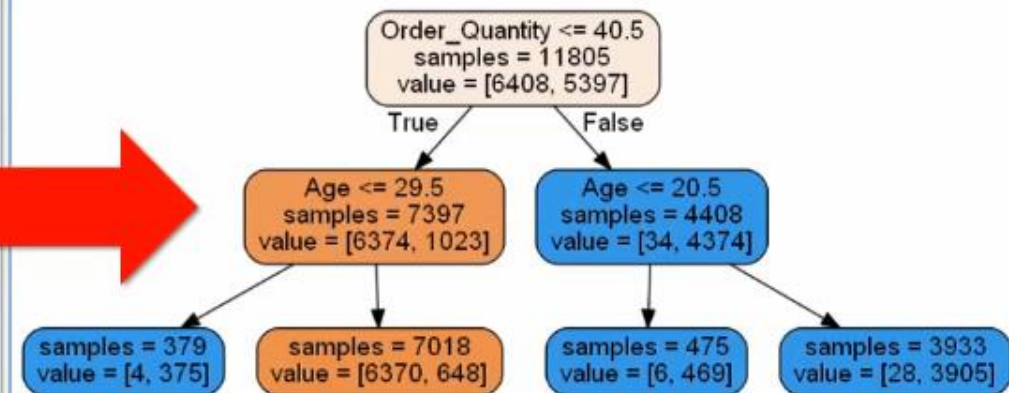
Console 1/A

```

...:
...: graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
...: Image(graph.create_png())
...:

```

Out[44]:



In [45]:

Python console History log IPython console


```

Spyder (Python 3.5)
File Edit Search Source Run Debug Consoles Projects Tools View Help
C:\Users\StatInfer
D:\Dropbox\0. Video Proj Files\Python ML\Session 3. Decision Tree\Active Presenter files\3.Decision_Trees_V3.py
3.Decision_Trees_V3.py*
1 #Import Data
2 import pandas as pd
3
4 ##Ecom_Cust_Survey = pd.read_csv('...',header = 0)
5 df = pd.read_csv('D:\\Google Drive\\Training\\Datasets\\Ecom_Cust_Relationship_Management\\Ecom_Cust_5
6
7 df.dropna(inplace='True') # to remove all the missing values rows..
8 #Q 1. How many customers have participated in the survey?
9 df.info()
10 #ANS: 11805
11
12 #Q.2 Overall most of the customers are satisfied or dis-satisfied?
13 #total number of customers
14 df.shape
15 df.head()
16
17
18 #number of satisfied customers
19 satisfied = df['Overall_Satisfaction'].map( {'Dis Satisfied': 0, 'Satisfied': 1} ).astype(int).sum()
20 satisfied
21 #number of dis-satisfied customers
22 df.shape[0]-satisfied
23
24
25 #Q 3. Can you segment the data and find the concentrated satisfied and dis-satisfied customer segments?
26 #solution:
27 # We will create a tree model in python using the sci-kit module
28 # before that we will need to convert most of the feature data into numerical or hash values as scikit
29 # Welcome to variable transformation
30
31 df['Region'] = df['Region'].map( {'EAST': 1, 'WEST': 2, 'NORTH': 3, 'SOUTH':4} ).astype(int)
32 df['Customer_Type'] = df['Customer_Type'].map({'Prime': 1, 'Non_Prime': 0}).astype(int)
33
34 #We will also need to change the column names, as '.' and spaces are part of many basic functions in py
35 df.rename(columns={'Order_Quantity':'Order_Quantity', 'Improvement_Area':'Improvement_Area'})
36 df['Improvement_Area'] = df['Improvement_Area'].map({'Improvement_Area':'Improvement_Area'})
37 df['Overall_Satisfaction'] = df['Overall_Satisfaction']
38
39 #Need the library to create the tree
40 from sklearn.tree import DecisionTreeClassifier
41
42 #Defining Features and Labels

```

Rule 1. Concentrated Satisfied

Name	Type	Size	Value
X	DataFrame	(11805, 4)	Column names: Age, Order_Quantity, Customer_Type, Improvement_Area
df	DataFrame	(11805, 7)	Column names: Cust_num, Region, Age, Order_Quantity, Customer_Type, Improvement_Area, Overall_Satisfaction
features	list	4	['Age', 'Order_Quantity', 'Customer_Type', 'Improvement_Area']
satisfied	int64	1	class 'numpy.int64'
y	Series	(11805,)	class 'pandas.core.series.Series'

Variable explorer

File explorer

Help

Python console

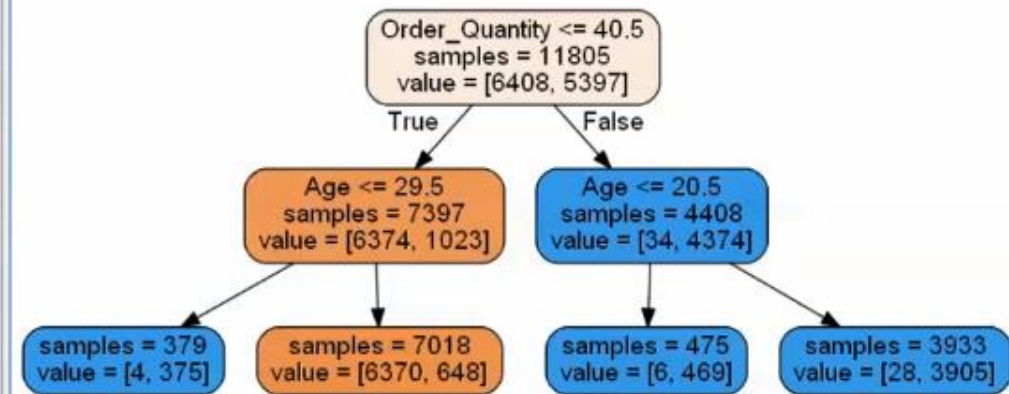
Console 1/A

```

...:
...: graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
...: Image(graph.create_png())
...:

```

Out[44]:



In [45]:

Python console History log Python console

Spyder (Python 3.5)

File Edit Search Source Run Debug Consoles Projects Tools View Help

Editor - D:\Dropbox\0_Video Proj Files\Python ML\Session 3. Decision Tree\Active Presenter files\3.Decision_Trees_V3.py

```
1 #Import Data
2 import pandas as pd
3
4 ##Ecom_Cust_Survey = pd.read_csv('...', header = 0)
5 df = pd.read_csv('D:\\Google Drive\\Training\\Datasets\\Ecom_Cust_Relationship_Management\\Ecom_Cust_5
6
7 df.dropna(inplace=True) # to remove all the missing values rows..
8 #Q 1. How many customers have participated in the survey?
9 df.info()
10 #ANS: 11805
11
12 #Q.2 Overall most of the customers are satisfied or dis-satisfied?
13 #total number of customers
14 df.shape
15 df.head()
16
17
18 #number of satisfied customers
19 satisfied = df['Overall_Satisfaction'].map( {'Dis Satisfied': 0, 'Satisfied': 1} ).astype(int).sum()
20 satisfied
21 #number of dis-satisfied customers
22 df.shape[0]-satisfied
23
24
25 #Q 3. Can you segment the data and find the concentrated satisfied and dis-satisfied customer segments?
26 #solution:
27 # We will create a tree model in python using the sci-kit module
28 # before that we will need to convert most of the feature data into numerical or hash values as scikit
29 # Welcome to variable transformation
30
31 df['Region'] = df['Region'].map( {'EAST': 1, 'WEST': 2, 'NORTH': 3, 'SOUTH':4} ).astype(int)
32 df['Customer_Type'] = df['Customer_Type'].map({'Prime': 1, 'Non_Prime': 0}).astype(int)
33
34 #We will also need to change the column names, as '.' and spaces are part of many basic functions in py
35 df.rename(columns={'Order Quantity':'Order_Quantity', 'Improvement Area': 'Improvement_Area'})
36 df['Improvement_Area'] = df['Improvement_Area'].map({'Website UI':1, 'Packaging':2, 'Product Quality':3})
37 df['Overall_Satisfaction'] = df['Overall_Satisfaction'].map( {'Dis Satisfied': 0, 'Satisfied': 1} ).astype(int)
38
39 #Need the library to create the tree
40 from sklearn.tree import DecisionTreeClassifier
41
42 #Defining Features and Labels
```

Variable explorer

Name	Type	Size	Value
X	DataFrame	(11805, 4)	Column names: Age, Order_Quantity, Customer_Type, Improvement_Area
df	DataFrame	(11805, 7)	Column names: Cust_num, Region, Age, Order_Quantity, Customer_Type, Improvement_Area, Overall_Satisfaction
features	list	4	['Age', 'Order_Quantity', 'Customer_Type', 'Improvement_Area']
satisfied	int64	1	class 'numpy.int64'
y	Series	(11805,)	class 'pandas.core.series.Series'

Variable explorer File explorer Help

Python console

Console 1/A

```
....:
....: graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
....: Image(graph.create_png())
....:
Out[44]:
```

```
graph TD
    Root["Order_Quantity <= 40.5  
samples = 11805  
value = [6408, 5397]"]
    Root -- True --> Node1["Age <= 29.5  
samples = 7397  
value = [6374, 1023]"]
    Root -- False --> Node2["Age <= 20.5  
samples = 4408  
value = [34, 4374]"]
    Node1 --> Leaf1["samples = 7018  
value = [6370, 648]"]
    Node2 --> Leaf2["samples = 475  
value = [6, 468]"]
    Node2 --> Leaf3["samples = 3933  
value = [28, 3905]"]
```

RULE 2 Concentrated Dissatisfied

In [45]:


```

1 #import Data
2 import pandas as pd
3
4 #Ecom_Cust_Survey = pd.read_csv('...', header = 0)
5 df = pd.read_csv('D:\\Google Drive\\Training\\Datasets\\Ecom_Cust_Relationship_Management\\Ecom_Cust_')
6
7 df.dropna(inplace=True) # to remove all the missing values rows..
8 #Q 1. How many customers have participated in the survey?
9 df.info()
10 INFO: 11805
11
12 #Q 2 Overall, most of the customers are satisfied or dis-satisfied?
13 #total number of customers
14 df.shape
15 df.head()
16
17
18 #number of satisfied customers
19 satisfied = df['Overall_Satisfaction'].map({'Dis Satisfied': 0, 'Satisfied': 1}).astype(int).sum()
20 satisfied
21 #number of dis-satisfied customers
22 df.shape[0]-satisfied
23
24
25 #Q 3. Can you segment the data and find the concentrated satisfied and dis-satisfied customer segments?
26 #solution:
27 # We will create a tree model in python using the scri-bit module
28 # before that we will need to convert most of the feature data into numerical or hash values as scri-bit
29 # welcome to variable transformation
30
31 df['Region'] = df['Region'].map({'EAST': 1, 'WEST': 2, 'NORTH': 3, 'SOUTH': 4}).astype(int)
32 df['Customer_Type'] = df['Customer_Type'].map({'Prime': 1, 'Non_Prime': 0}).astype(int)
33
34 #we will also need to change the column names, as '.' and spaces are part of many basic functions in p
35 df.rename(columns={'Order_Quantity': 'Order_Quantity', 'Improvement_Area': 'Improvement_Area'}, inplace
36 df['Improvement_Area'] = df['Improvement_Area'].map({'Website UI': 1, 'Packing & Shipping': 2, 'Product
37 df['Overall_Satisfaction'] = df['Overall_Satisfaction'].map({'Dis Satisfied': 0, 'Satisfied': 1}).an
38
39 #need the library to create the tree
40 from sklearn.tree import DecisionTreeClassifier
41
42 #Defining Features and Labels

```

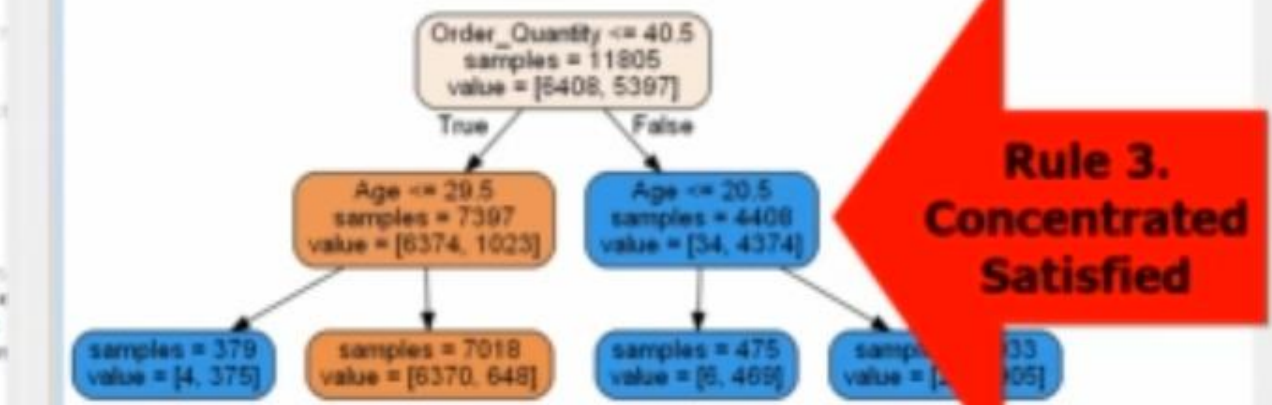
Name	Type	Size	Value
X	DataFrame	(11805, 4)	Column names: Age, Order_Quantity, Customer_Type, Improvement_Area
df	DataFrame	(11805, 7)	Column names: Cust_num, Region, Age, Order_Quantity, Customer_Type, Improvement_Area, Overall_Satisfaction
features	list	4	['Age', 'Order_Quantity', 'Customer_Type', 'Improvement_Area']
satisfied	int64	1	class 'numpy.int64'
y	Series	(11805,)	class 'pandas.core.series.Series'

```

----
.... graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
.... image(graph.create_png())
----

```

Out[44]:



In [45]:

```

19 satisfied = df['Overall_Satisfaction'].map( {'Dis Satisfied': 0, 'Satisfied': 1} ).astype(int).sum()
20 satisfied
21 #number of dis-satisfied customers
22 df.shape[0]-satisfied
23
24
25 #Q 3. Can you segment the data and find the concentrated satisfied and dis-satisfied customer segment:
26 #solution:
27 # We will create a tree model in python using the sci-kit module
28 # before that we will need to convert most of the feature data into numerical or hash values as scikit
29 # Welcome to variable transformation
30
31 df['Region'] = df['Region'].map( {'EAST': 1, 'WEST': 2, 'NORTH': 3, 'SOUTH':4} ).astype(int)
32 df['Customer_Type'] = df['Customer_Type'].map({'Prime': 1, 'Non_Prime': 0}).astype(int)
33
34 #We will also need to change the column names, as '.' and spaces are part of many basic functions in py
35 df.rename(columns={'Order_Quantity':'Order_Quantity', 'Improvement Area ':'Improvement_Area'}, inplace
36 df['Improvement_Area'] = df['Improvement_Area'].map({'Website UI':1, 'Packing & Shipping':2, 'Product
37 df['Overall_Satisfaction'] = df['Overall_Satisfaction'].map( {'Dis Satisfied': 0, 'Satisfied': 1} ).as
38
39 #Need the library to create the tree
40 from sklearn.tree import DecisionTreeClassifier
41
42 #Defining Features and Labels
43 features= list(df.columns[2:6])
44 features
45
46 X = df[features]
47 y = df['Overall_Satisfaction']
48
49 #Building Tree Model
50 clf = DecisionTreeClassifier(max_depth=2)
51 clf.fit(X,y)
52
53 #What are the major characteristics of satisfied customers?
54
55 #Plotting the trees
56 #Unfortunately drawing a beautiful tree is not easy in python, Still
57 #you will need to install pydot
58 #use this command in your anaconda prompt: conda install -c anaconda pydot=1.0.28
59
60 from IPython.display import Image

```

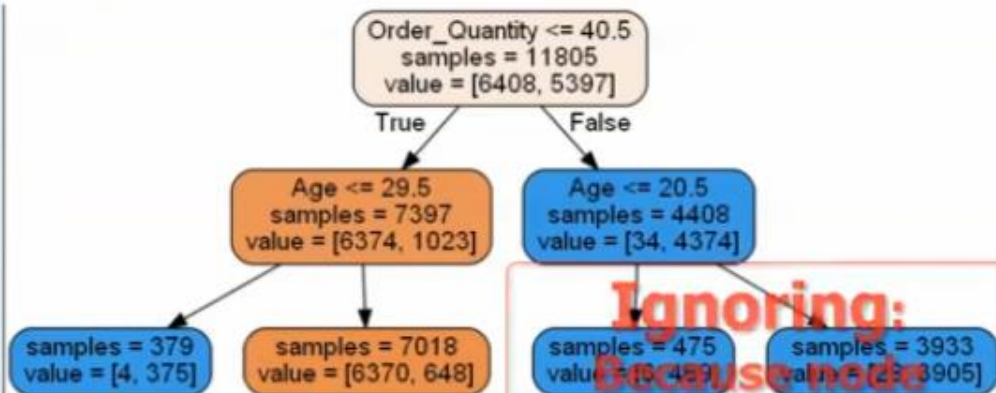
Name	Type	Size	Value
X	DataFrame	(11805, 4)	Column names: Age, Order_Quantity, Customer_Type, Improvement_Area
df	DataFrame	(11805, 7)	Column names: Cust_num, Region, Age, Order_Quantity, Customer_Type, Improvement_Area, Overall_Satisfaction
features	list	4	['Age', 'Order_Quantity', 'Customer_Type', 'Improvement_Area']
satisfied	int64	1	class 'numpy.int64'
y	Series	(11805,)	class 'pandas.core.series.Series'

```

...:
...: graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
...: Image(graph.create_png())
...:

```

Out[44]:

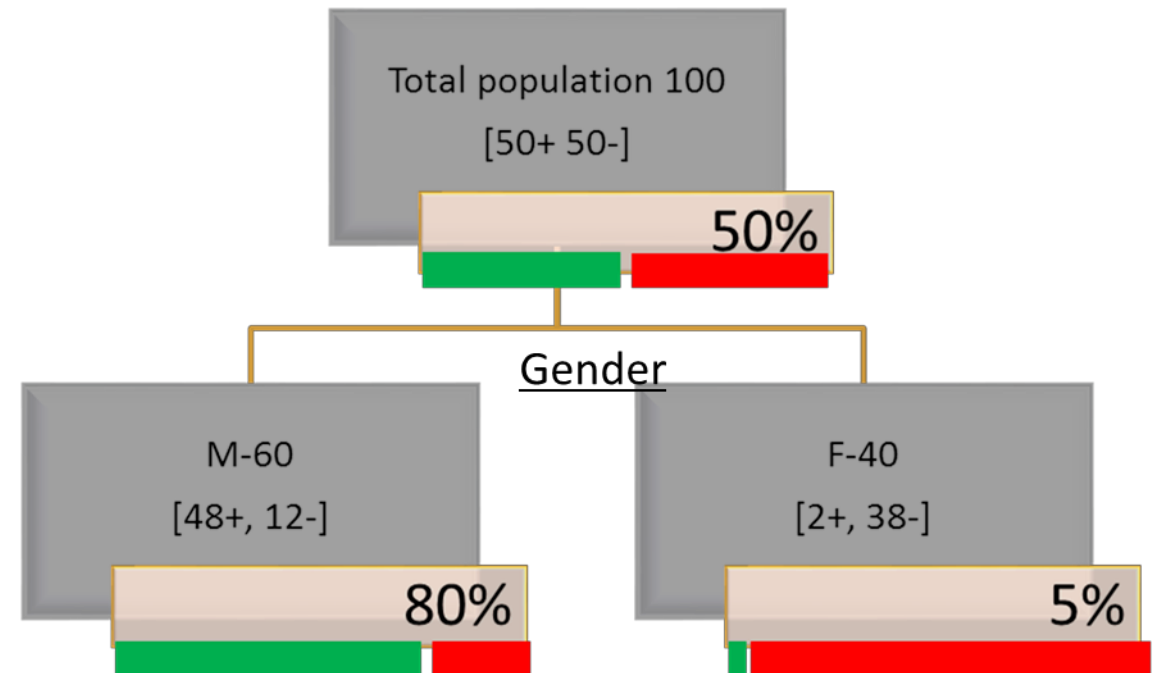


In [45]:

LAB: Entropy Calculation - use calculator or Excel

LAB Entropy Calculation

- Calculate entropy at the root for the given population
- Calculate the entropy for the two distinct gender segments



Code- Entropy Calculation

- Entropy at root 100%
- Male Segment : $(-48/60)*\log(48/60,2)-(12/60)*\log(12/60,2)$
 - 0.7219281
- FemaleSegment : $(-2/40)*\log(2/40,2)-(38/40)*\log(38/40,2)$
 - 0.286397

Thank you
